



# **Interfacing the Storage Resource Broker (SRB) to the Hierarchical Resource Manager (HRM)**

**Arie Shoshani, Alex Sim (LBNL)  
Reagan Moore, Bing Zhu (SDSC)**

**PPDG meeting  
July 13-14, 2000**

PPDG meeting, July 2000

## **Outline**



- **HRM role in the Data Grid architecture**
- **HRM API description**
- **HRM - HPSS system functionality**
- **SRB Client calls**
- **SRB-HRM architecture**
- **Status**

PPDG meeting, July 2000

## HRM role in the Data Grid architecture



- The class of Storage Resource Managers (SRMs) includes:
  - HRM: for managing the access to tape resources
    - may or may not have a disk cache
    - functionality generic but needs to be specialize for specific mass storage systems
    - e.g. HRM-HPSS, HRM-Enstore, ...
  - DRM: for managing disk resources
    - functionality generic but needs to be specialize for specific disk systems
    - e.g. DRM-FileSystem, DRM-DPSS, ...

PPDG meeting, July 2000

## HRM role in the Data Grid architecture



### Functionality Examples

#### HRM functionality may include :

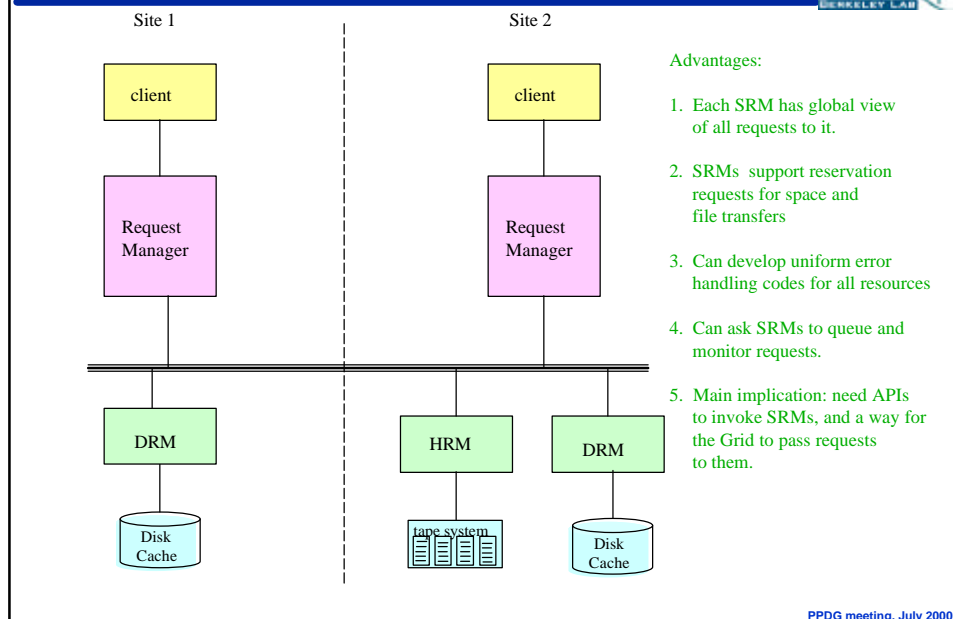
- queuing of file transfer requests (kind of a reservation)
- reordering of requests to optimize PFTP order (ordered by files on the same tape)
- Monitoring progress and error messages
- rescheduling transfers that failed

#### DRM functionality may include:

- keeping tracks of files in cache
- managing space reservations
- making decisions on which files to remove when space is needed
- optimizing cache use - sharing files requested by multiple clients
- enforce local policy for cache use

PPDG meeting, July 2000

## HRMs & DRMs in a Grid Architecture



PPDG meeting, July 2000

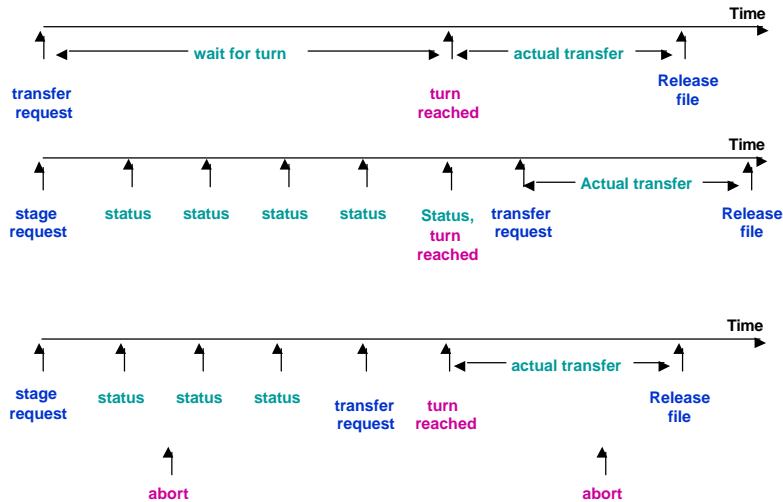
## HRM API description



- **API Functionality**
  - Request to transfer a file to destination disk
    - A blocking call
  - Request to stage a file to HRM disk
    - A non-blocking call when HRM disk exists
  - Request status/time estimate
    - How long before file request will be processed
  - Request to abort a file transfer or stage
    - In case file no longer needed
  - Release a file
    - After file was moved to destination
    - optional to improve system efficiency
  - Call back when file is staged

PPDG meeting, July 2000

## Examples of call sequences



PPDG meeting, July 2000

## CORBA IDL for HRM



```
// hrm.idl
// IDL for HPSS Resource Manager
//
#ifdef HRM_IDL
#define HRM_IDL

enum hrmStatus {
    HRM_TRANSFER_COMPLETED,
    HRM_TRANSFER_FAILED,
    HRM_STAGE_ACCEPTED,
    HRM_STAGE_COMPLETED,
    HRM_STAGE_REFUSED,
    HRM_STAGE_FAILED,
    HRM_RELEASE_DONE,
    HRM_RELEASE_FAILED,
    HRM_ABORT_DONE,
    HRM_ABORT_FAILED,
    HRM_FILE_DOES_NOT_EXIST,
    HRM_HPSS_DOWN,
    HRM_HPSS_ERROR
};

module HPSSResourceManager {
    interface hrmGrid {
        hrmStatus transferFile
            (in string sourceUrl,
             in string destinationUrl,
             in long waitTime,
             in short userPriority,
             out double timeEstimate);
        hrmStatus stageFile
            (in string reference,
             in string destinationUrl,
             in long waitTime,
             in short userPriority,
             out double timeEstimate);
        boolean abort(in string destinationUrl);
        double timeEstimate(in string destinationUrl);
        boolean releaseFile(in string destinationUrl);
    };
};
#endif
```

PPDG meeting, July 2000

## HRM-HPSS system functionality



- All transfers go through HRM disk
  - reasons: flexibility of pre-staging
  - disk is sufficiently cheap for a large cache
  - opportunity to optimize for same file requests
- Functionality
  - queuing file transfers
  - file queue management
  - File clustering parameter
  - Transfer rate estimation
  - Query estimation - total time
  - Error handling

PPDG meeting, July 2000

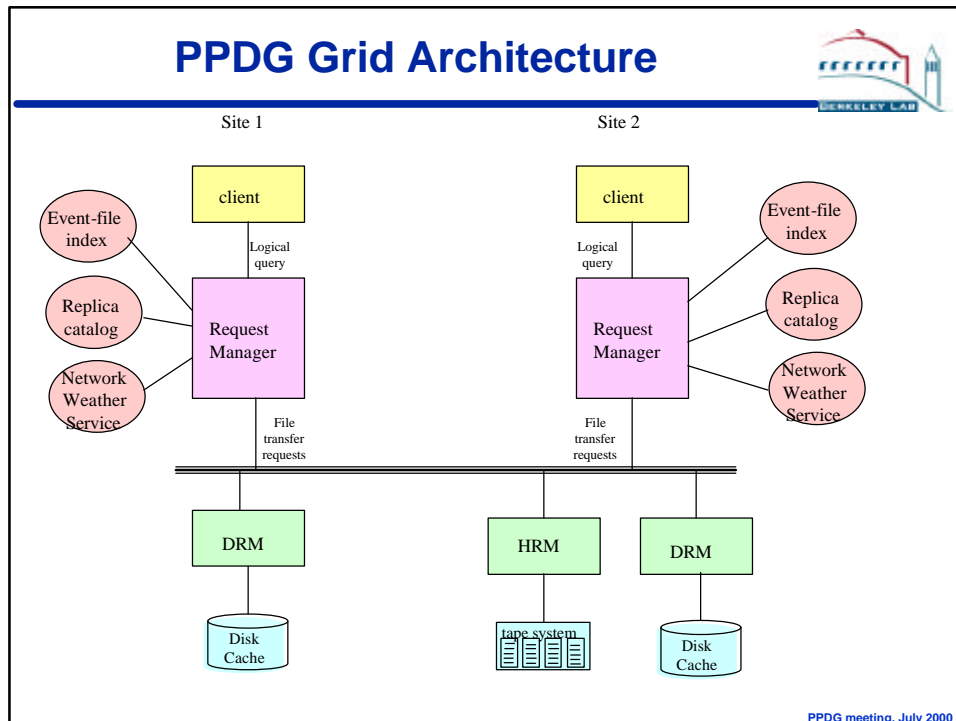
## SRB Client calls



- Sget
  - blocking call (wait till file is transferred)
  - should only be issued if space on client disk is allocated
- Sstage
  - non-blocking call (returns time estimate)
  - requires a server if want to be notified
  - can issue status to find out if file was cached
- Sstatus
  - returns time estimate for file to be staged to HRM's disk
- Sabort
  - cancel this file request

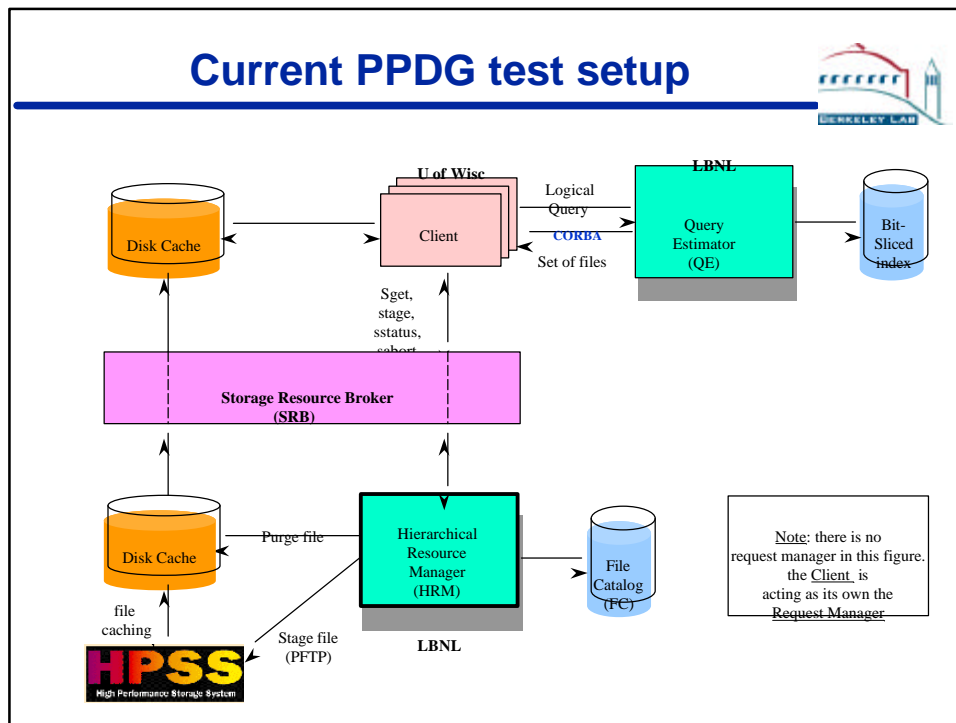
PPDG meeting, July 2000

## PPDG Grid Architecture



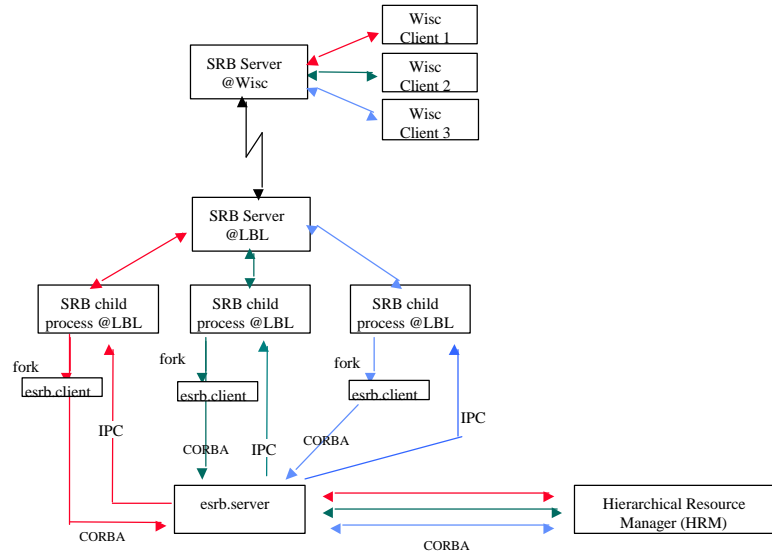
PPDG meeting, July 2000

## Current PPDG test setup

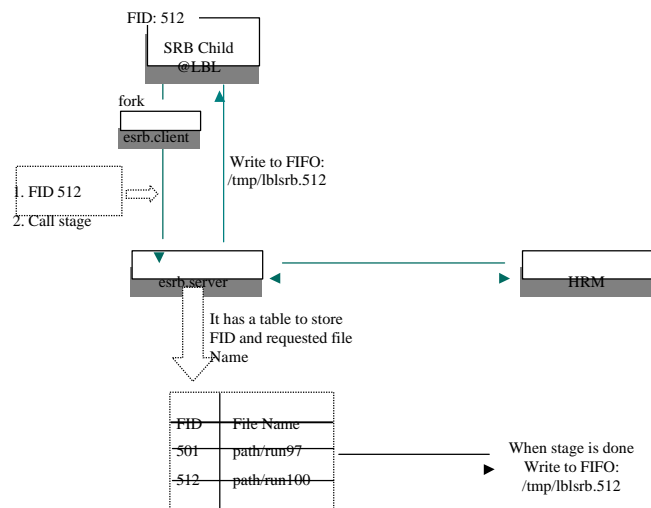


PPDG meeting, July 2000

## SRB-HRM interaction



PPDG meeting, July 2000



PPDG meeting, July 2000

## Status



- **Functionality that was tested**
  - request a file transfer - Sget
  - A single Client
- **Functionality that exists but was not tested**
  - Sget, Sstage, Sstatus
  - multiple Clients concurrent requests
  - use of replica catalog developed by SDSC
- **Functionality that does not exist**
  - A request manager that uses the replica catalog
  - A request manager that accesses multiple sites
  - A request manager that makes informed choices
  - A disk resource manager

PPDG meeting, July 2000

## Optional Slides



- The following slides shows how HRM implements the functionality it provides
- They will be shown only if time permits

PPDG meeting, July 2000



## Queuing File Transfers



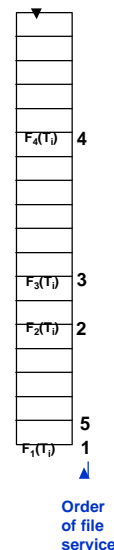
- Number of PFTPs to HPSS are limited
  - limit set by a parameter - No\_PFTP
  - parameter can be changed dynamically
- HRM is multi-threaded
  - issues and monitors multiple PFTPs in parallel
- All requests beyond PFTP limit are queued
- File Catalog used to provide for each file
  - HPSS path/file\_name
  - Disk cache path/file\_name
  - File size
  - tape ID

PPDG meeting, July 2000

## File Queue Management

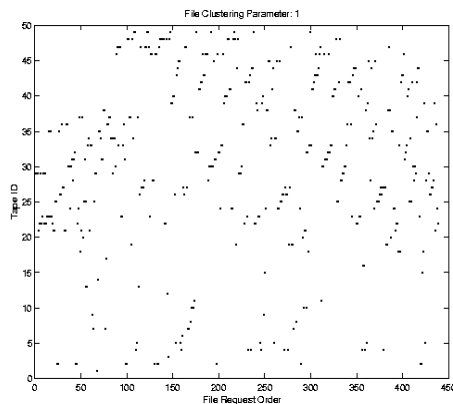


- Goal
  - minimize tape mounts
  - still respect the order of requests
  - do not postpone unpopular tapes forever
- File clustering parameter - FCP
  - If the file at top of queue is in Tape<sub>i</sub> and FCP > 1 (e.g. 4) then up to 4 files from Tape<sub>i</sub> will be selected to be transferred next
  - then, go back to file at top of queue
- Parameter can be set dynamically

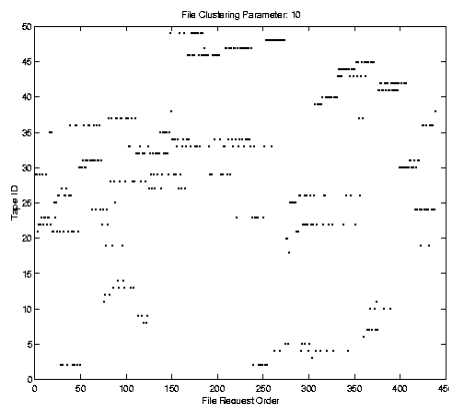


PPDG meeting, July 2000

## File Caching Order for different File Clustering Parameters



File Clustering Parameter = 1



File Clustering Parameter = 10

PPDG meeting, July 2000

## Transfer Rate ( $Tr$ ) Estimates



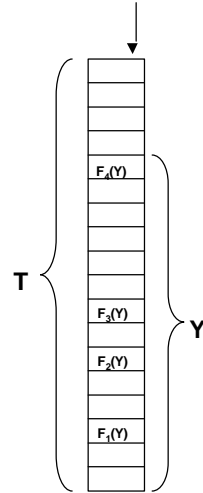
- Need  $Tr$  to estimate total time of a query
- $Tr$  is average over recent file transfers from the time PFTP request is made to the time transfer completes. This includes:
  - mount time, seek time, read to HPSS Raid, transfer to local cache over network
- For dynamic network speed estimate
  - check total bytes for all file being transferred over small intervals (e.g. 15 sec)
  - calculate moving average over  $n$  intervals (e.g. 10 intervals)

PPDG meeting, July 2000

## Query Estimate

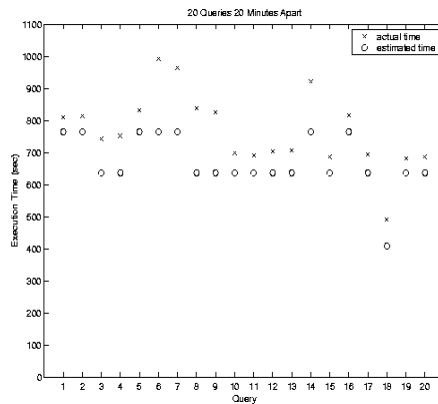


- Given: transfer rate  $Tr$ .
- Given a query for which:
  - $X$  files are in cache
  - $Y$  files are in the queue
  - $Z$  files are not scheduled yet
- Let  $s(\text{file\_set})$  be the total byte size of all files in  $\text{file\_set}$
- If  $Z = 0$ , then
  - $QuEst = s(Y)/Tr$
- If  $Z \neq 0$ , then
  - $QuEst = (s(T) + q.s(Z))/Tr$   
where  $q$  is the number of active queries



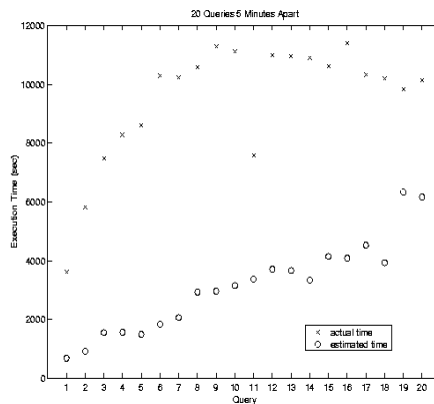
PPDG meeting, July 2000

## Reason for $q.s(Z)$



20 Queries of length ~20 minutes  
launched 20 minutes apart

Estimate pretty close



20 Queries of length ~20 minutes  
launched 5 minutes apart

Estimate bad - request  
accumulate in queue

PPDG meeting, July 2000

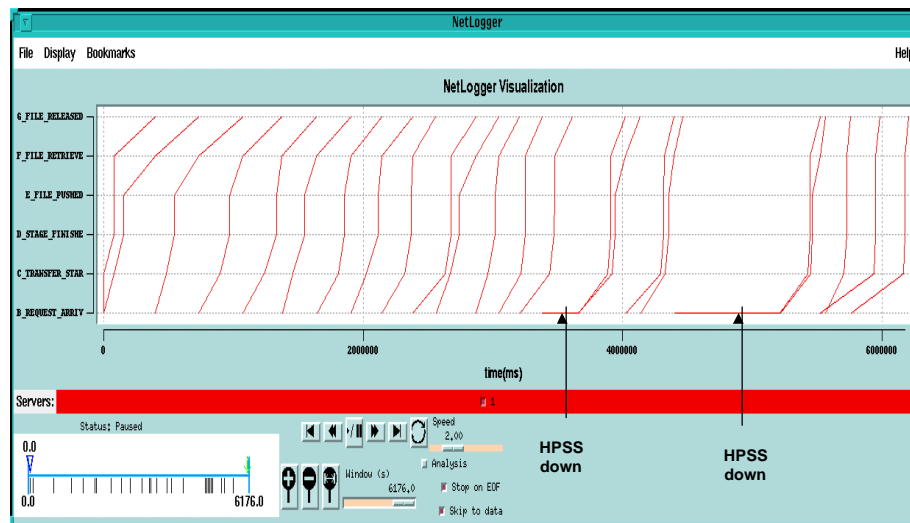
## Error Handling



- **5 generic errors**
  - **file not found**
    - return error to caller
  - **limit PFTP reached**
    - can't login
    - re-queue request, try later (1-2 min)
  - **HPSS error (I/O, device busy)**
    - remove part of file from cache, re-queue
    - try n times (e.g. 3), then return error "transfer\_failed"
  - **HPSS down**
    - re-queue request, try repeatedly till successful
    - respond to File\_status request with "HPSS\_down"

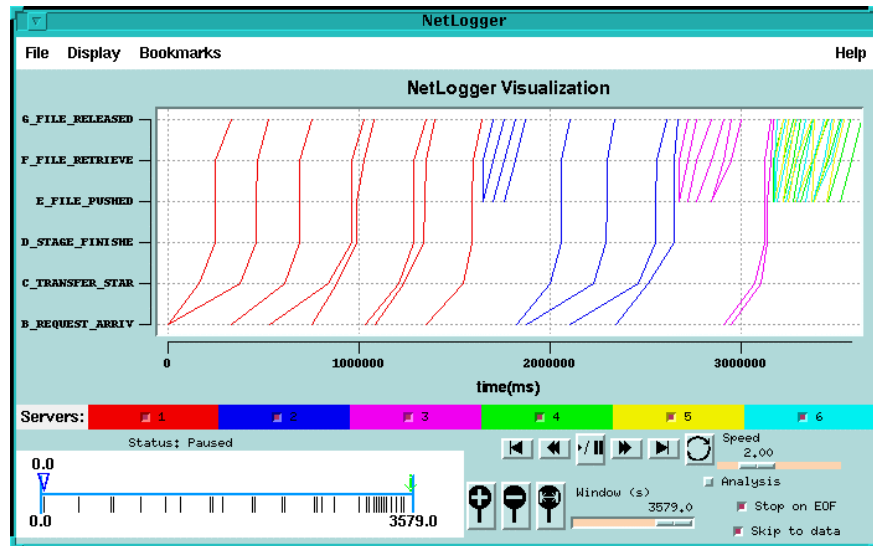
PPDG meeting, July 2000

## File Tracking (1)



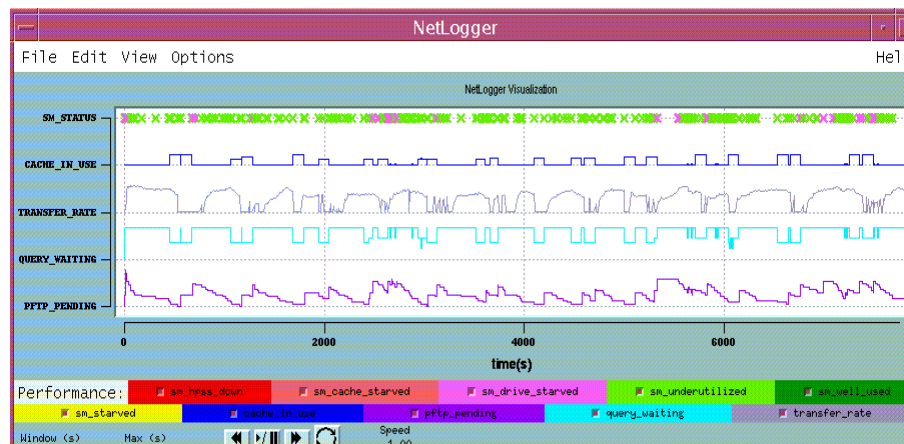
PPDG meeting, July 2000

## File Tracking (2)



PPDG meeting, July 2000

## Dynamic Display of Various Measurements



PPDG meeting, July 2000